



Segmentez les clients d'un site e-commerce
Synthèse

Contexte & Objectifs

PROJET

Sur la base des transactions effectuées sur 2 années pleines (septembre 2016 à août 2018), établir une **segmentation** des clients permettant de décrire leur comportement d'achat à des fins de marketing. Le nombre de classes devra être suffisamment élevé pour permettre une **stratégie diversifiée** en fonction des segments, tout en n'étant pas excessif pour rester **interprétable** dans le contexte métier.

DÉMARCHE

Après une analyse exploratoire des données, plusieurs modèles de **classification non supervisée** (clustering) seront testés et évalués en fonction de leur capacité à répondre à la problématique métier.



Sommaire

1 – ANALYSE EXPLORATOIRE DES DONNÉES

- 1.1 – EXPLORATION & NETTOYAGE
- 1.2 – FEATURE ENGINEERING & SÉLECTION
- 1.3 – STATISTIQUES DESCRIPTIVES

2 – MODÉLISATION

- 2.1 – PROBLÉMATIQUE & MÉTHODOLOGIE
- 2.2 – CRITÈRES MÉTIER & CHOIX DU MODÈLE FINAL
- 2.3 – CARACTÉRISTIQUES RFM DES CLUSTERS & STRATÉGIE MARKETING

3 – ÉVALUATION DU MODÈLE FINAL

- 3.1 – STABILITÉ DES CLUSTERS À L'INITIALISATION
- 3.2 – CONTRAT DE MAINTENANCE – MÉTHODOLOGIE & RECOMMANDATIONS

CONCLUSION & PERSPECTIVES

ANNEXES

1 – ANALYSE EXPLORATOIRE DES DONNÉES

1.1 – EXPLORATION & NETTOYAGE

JEU DE DONNEES INITIAL* :

5 tables de dimensions



Données géographiques – table ‘geolocation’

1,000,163 lignes x 6 colonnes – 0% NaN



Clients – table ‘customers’

99,441 lignes x 2 colonnes – 0% NaN



Vendeurs – table ‘sellers’

3,95 lignes x 2 colonnes – 0% NaN



Description produits – table ‘products’

32,951 lignes x 9 colonnes – 0.8% NaN



Catégories produits – table ‘categories’

71 lignes x 2 colonnes – 0% NaN

4 tables de faits



Ventes – table ‘baskets’

112,650 lignes x 7 colonnes – 0% NaN



Paielements – table ‘payments’

103,886 lignes x 5 colonnes – 0% NaN



Revue de produit – table ‘reviews’

99,224 lignes x 7 colonnes – 21% NaN



Livraisons – table ‘shipping’

99,441 lignes x 8 colonnes – 0.6% NaN

1 – ANALYSE EXPLORATOIRE DES DONNÉES

1.2 – FEATURE ENGINEERING & SÉLECTION

FEATURE ENGINEERING

- Longueur des titres et reviews
- Distance entre clients et vendeurs
- Groupement des produits en 13 puis 3 catégories
- Labels pour les transactions faites le weekend ou le soir après 18h
- Nombre de jours depuis la dernière transaction du client
- Ajout de données externes etc...*

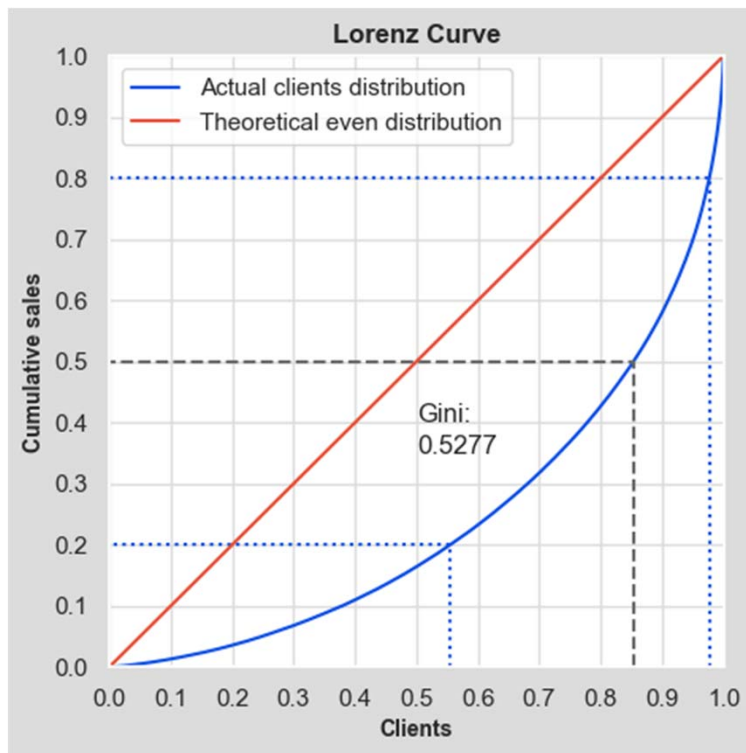
RETRAITEMENTS

- Suppressions des valeurs manquantes & doublons
- Suppression des variables trop corrélées entre elles
- Jointures entre les tables
- Agrégation des données par client
- Création itérative de plusieurs jeux de données avec différentes features numériques & catégorielles

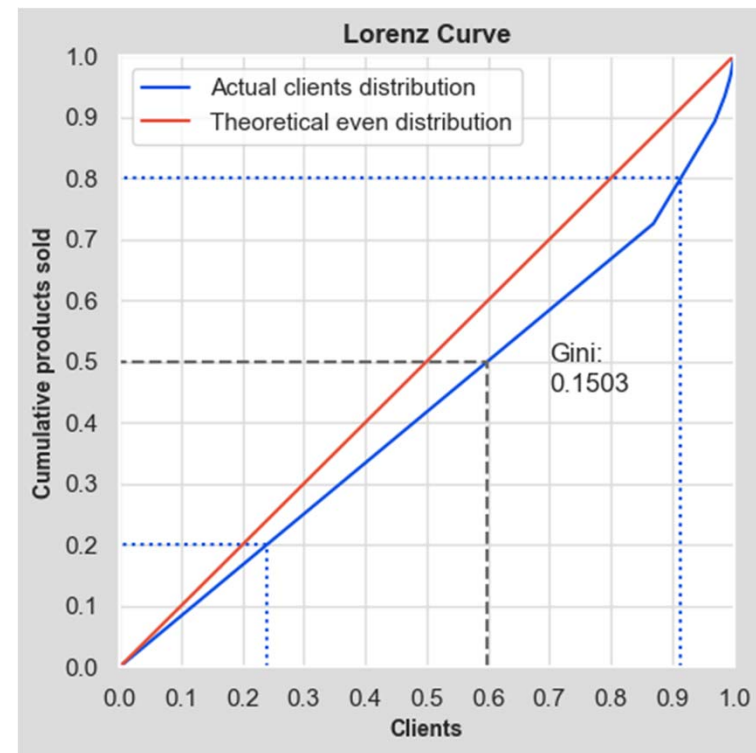
1 - ANALYSE EXPLORATOIRE DES DONNÉES

1.3 - STATISTIQUES DESCRIPTIVES - CHIFFRE D'AFFAIRES EN VALEUR & VOLUME

Valeur - Chiffre d'affaires total :

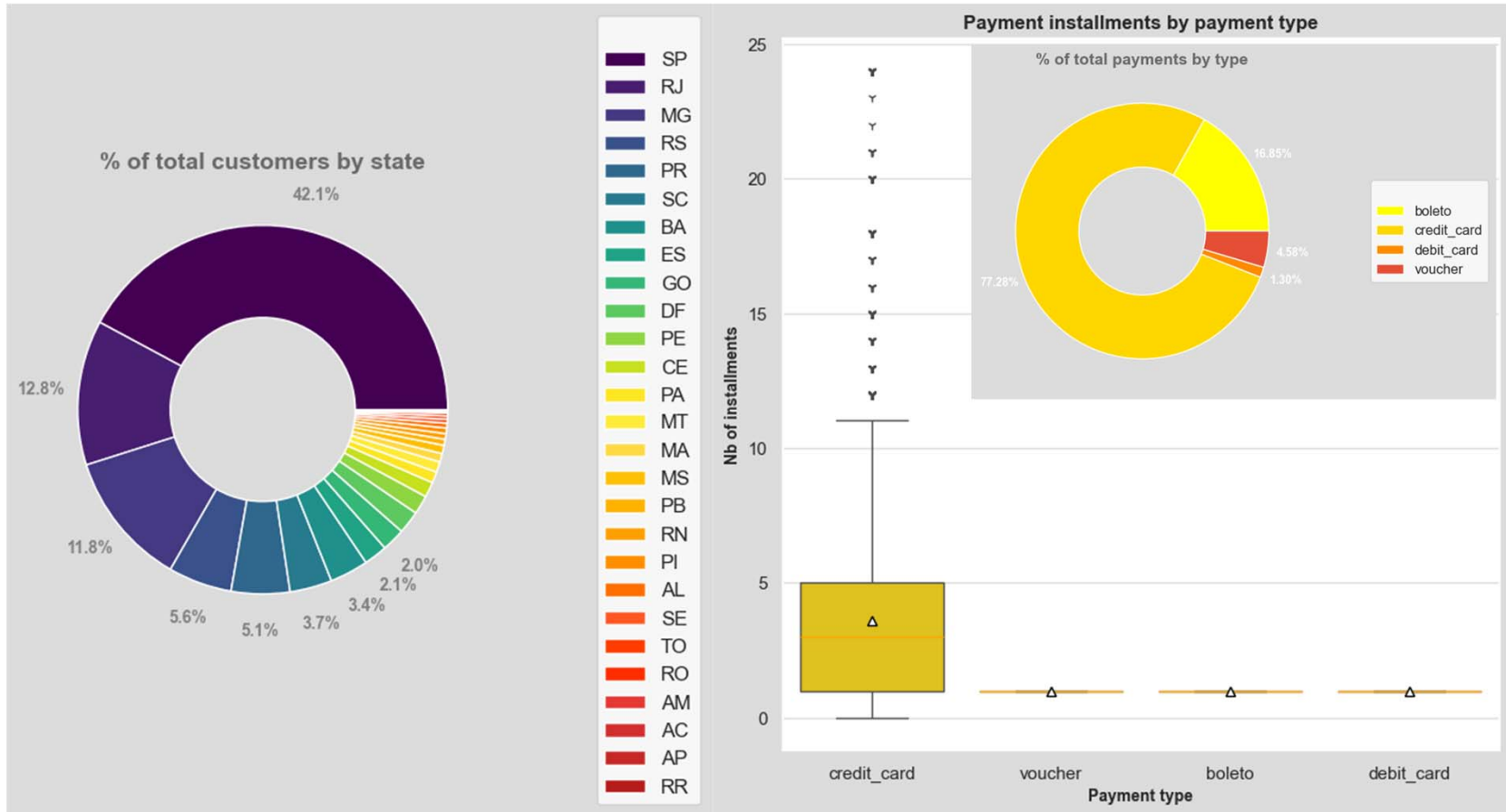


Volume - Nombre de produits vendus :



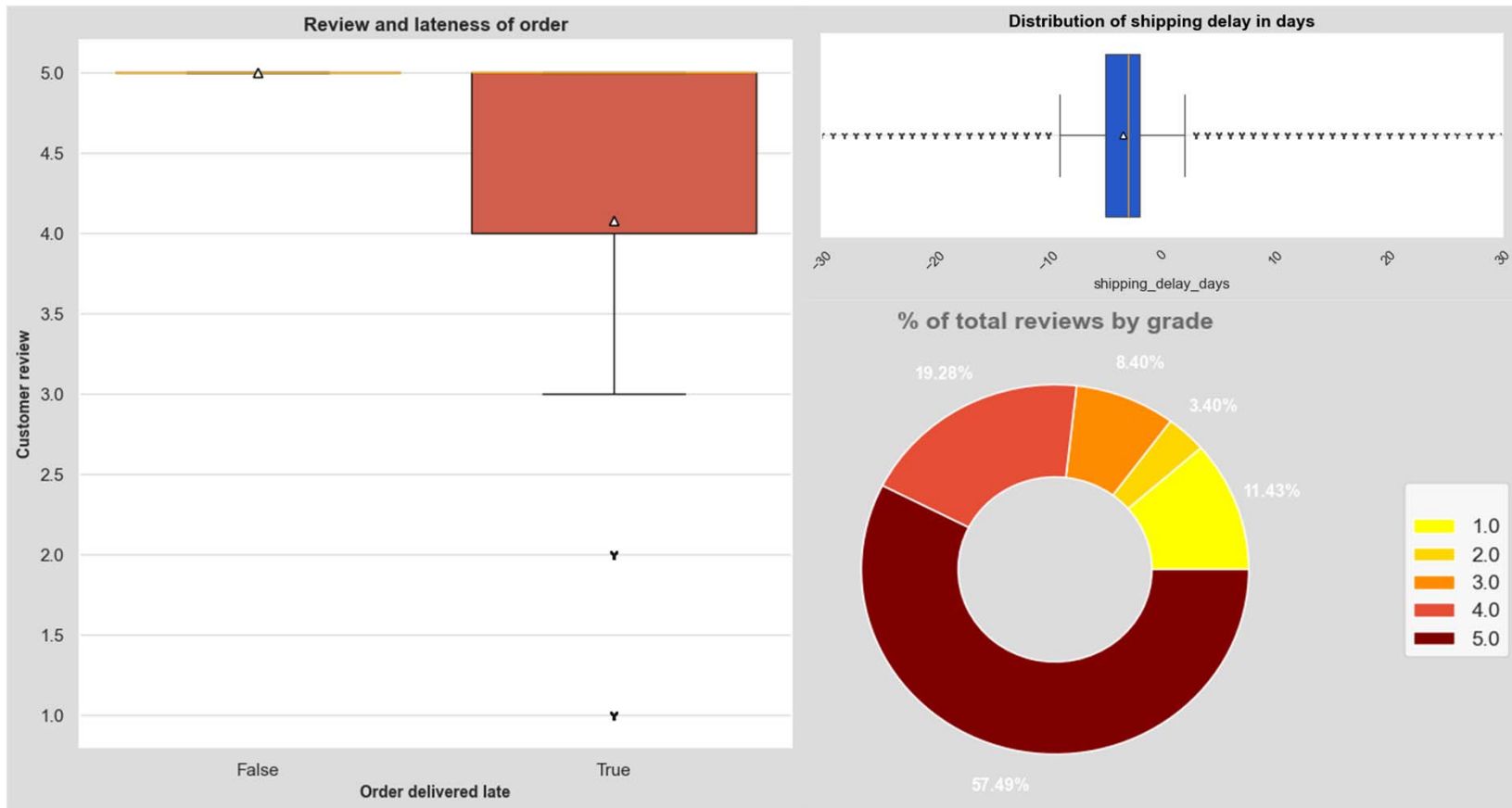
1 - ANALYSE EXPLORATOIRE DES DONNÉES

1.3 - STATISTIQUES DESCRIPTIVES - LOCALISATION & TYPES DE PAIEMENTS



1 - ANALYSE EXPLORATOIRE DES DONNÉES

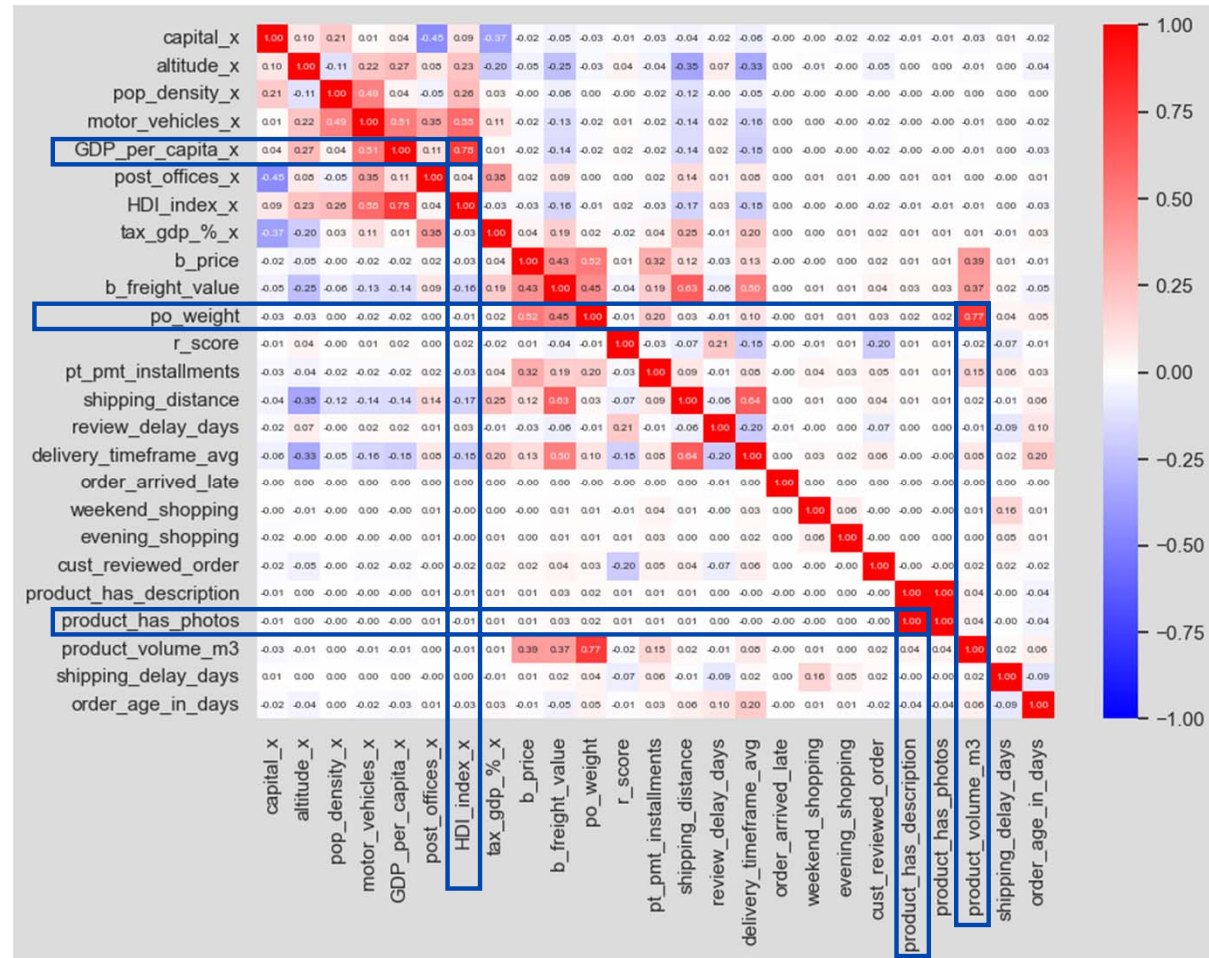
1.3 - STATISTIQUES DESCRIPTIVES - REVIEWS & RETARDS D'ACHEMINEMENT



1 - ANALYSE EXPLORATOIRE DES DONNÉES

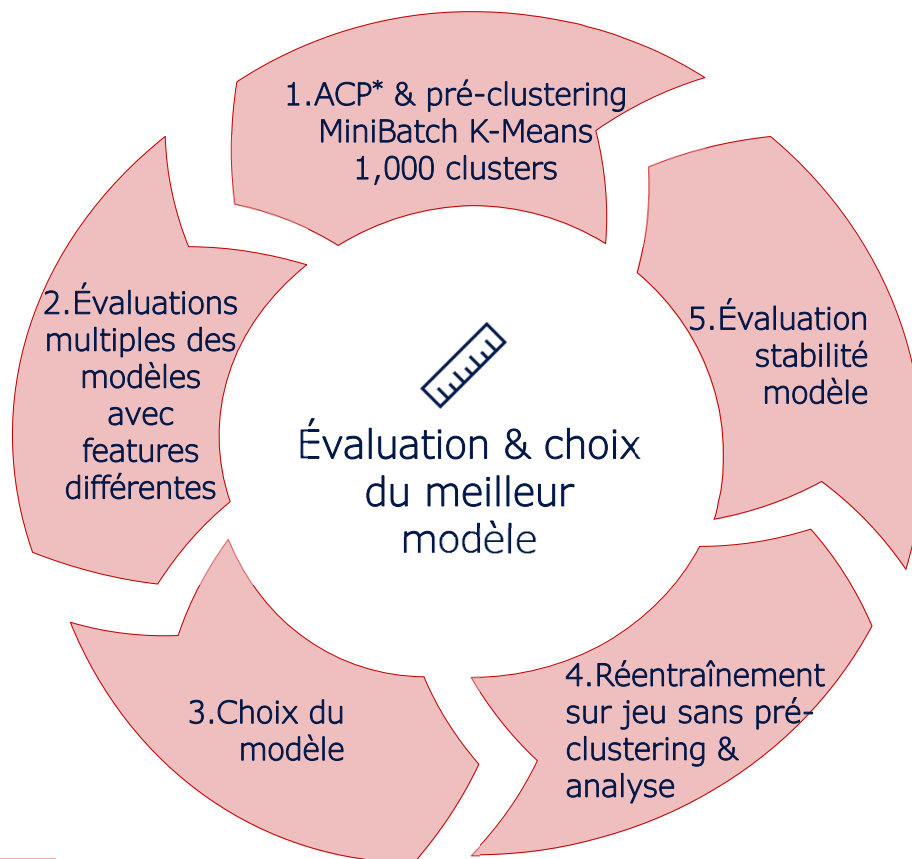
1.3 - STATISTIQUES DESCRIPTIVES - CORRÉLATIONS

- ❖ Suppression des features fortement corrélées
- ❖ Skew
- ❖ Variables non-gaussiennes
 - Coefficient de Spearman



2 – MODÉLISATION

2.1 – PROBLÉMATIQUE & MÉTHODOLOGIE



2. MÉTRIQUES D'ÉVALUATION DU NOMBRE DE CLUSTERS OPTIMAL:

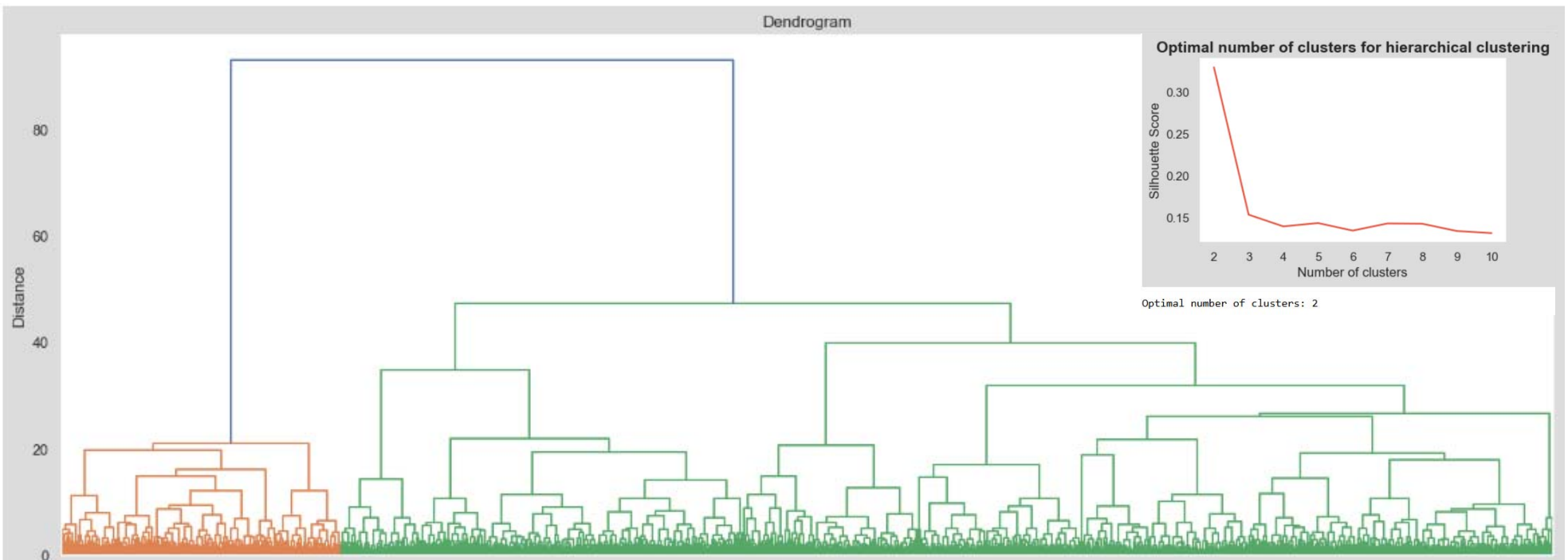
- Méthode du 'coude'
- Coefficient Silhouette
- Indice de Davies-Bouldin
- Indice de Kalinski-Harabasz

3. MODÈLES TESTÉS:

- Classification Ascendante Hiérarchique
- DBSCAN
- HDBSCAN
- Affinity Propagation
- K-Prototypes
- K-means

2 – MODÉLISATION

2.2 – CRITÈRES MÉTIER & CHOIX DU MODÈLE FINAL – CLASSIFICATION ASCENDANTE HIÉRARCHIQUE



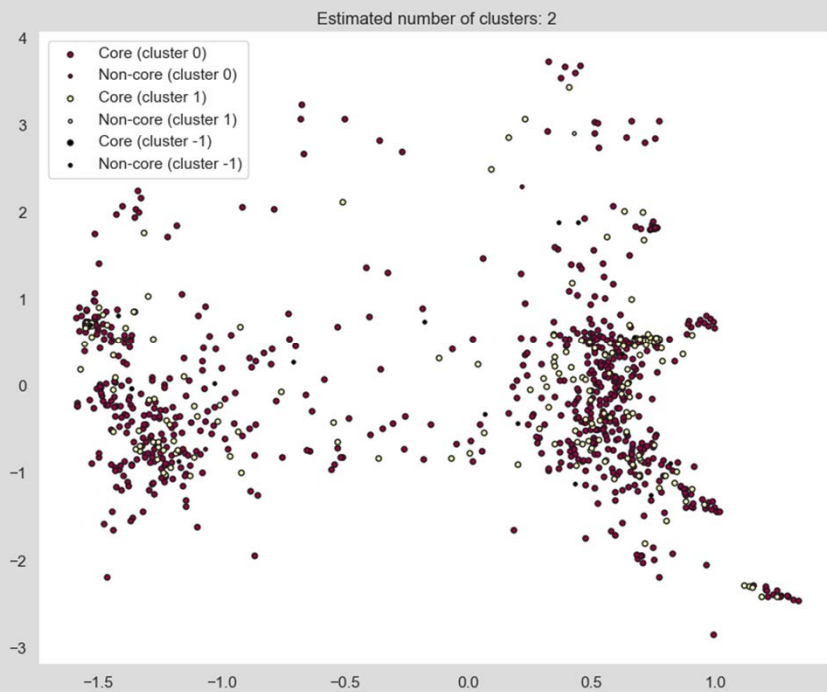
Nb of customer pre-clusters in 1st cluster : 187 . Population : 12437

Nb of customer pre-clusters in 2nd cluster : 813 . Population : 82903

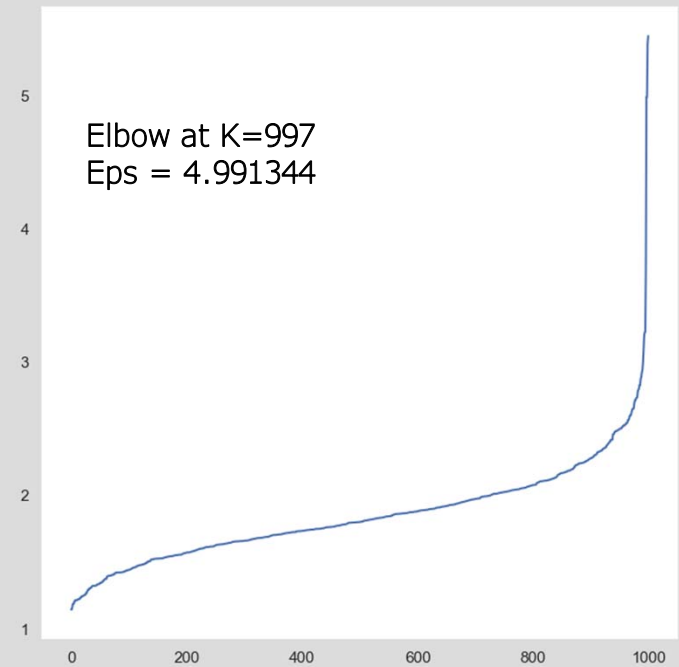
Total customers : 95340 . Total pre-clusters : 1000

2 – MODÉLISATION

2.2 – CRITÈRES MÉTIER & CHOIX DU MODÈLE FINAL - DBSCAN

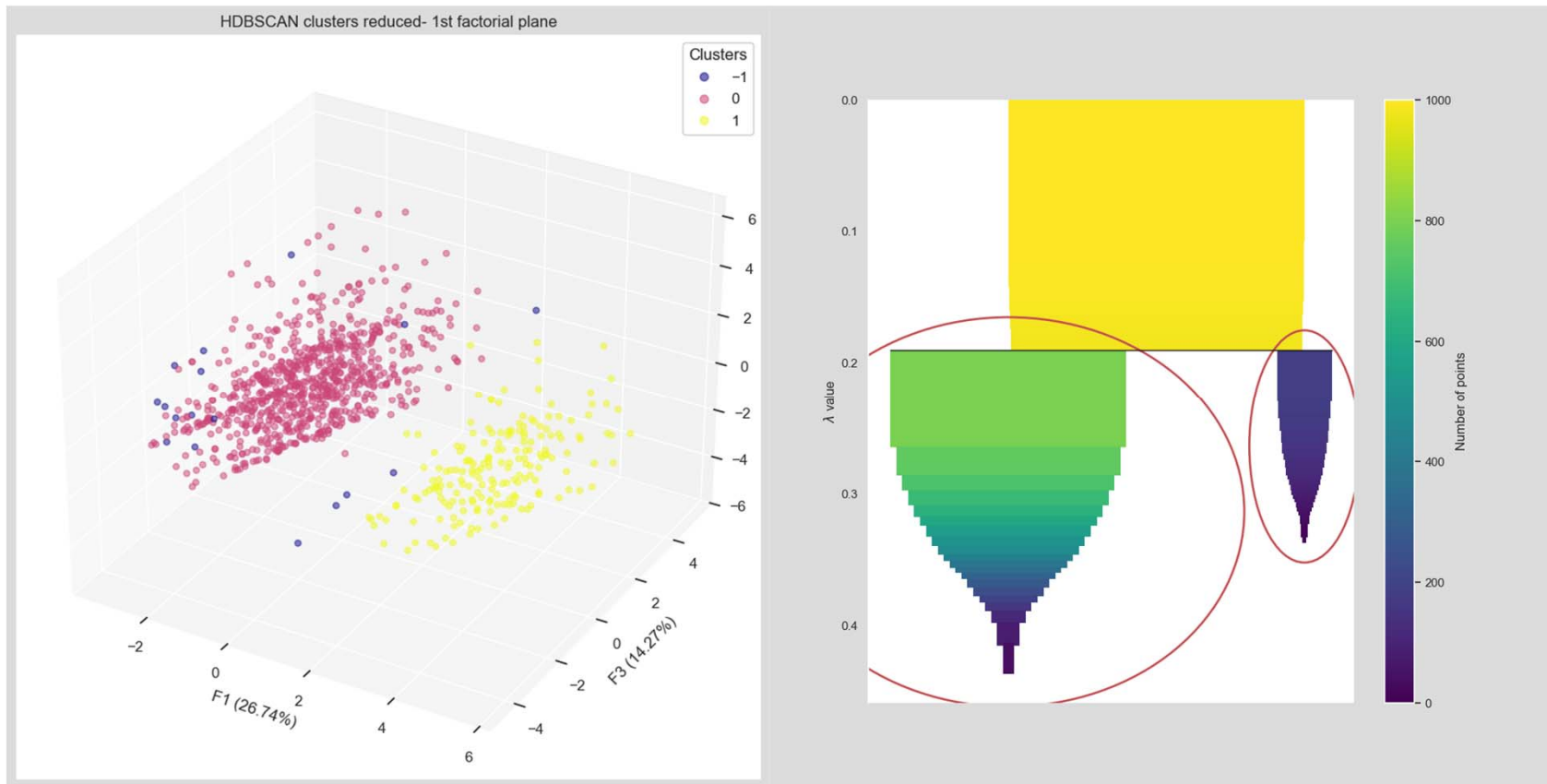


DBSCAN EPS & MinSamples



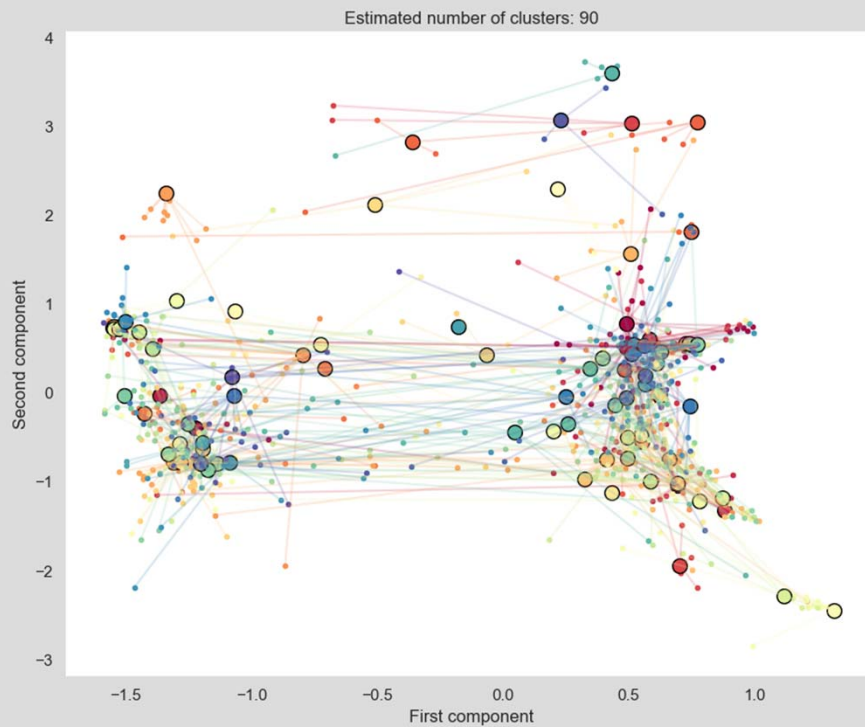
2 – MODÉLISATION

2.2 – CRITÈRES MÉTIER & CHOIX DU MODÈLE FINAL - HDBSCAN

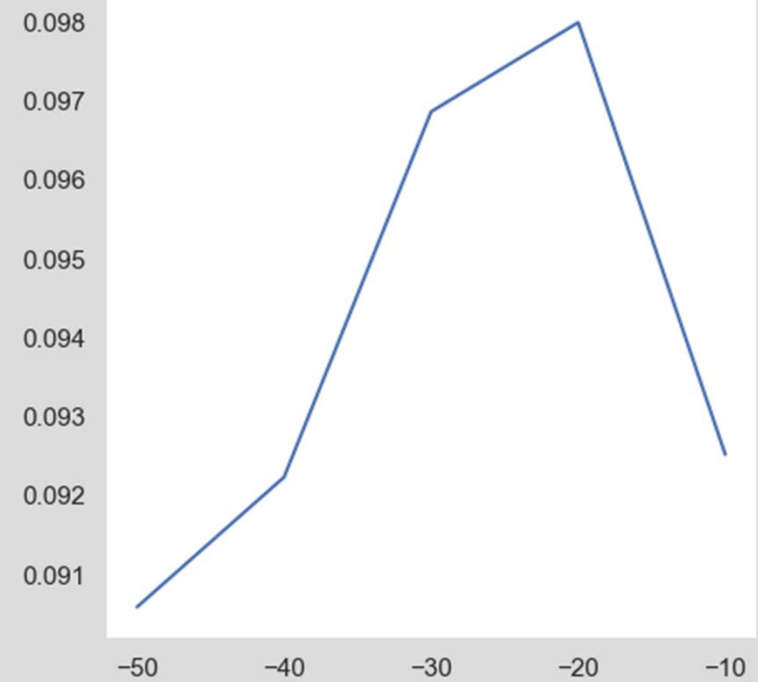


2 – MODÉLISATION

2.2 – CRITÈRES MÉTIER & CHOIX DU MODÈLE FINAL – AFFINITY PROPAGATION

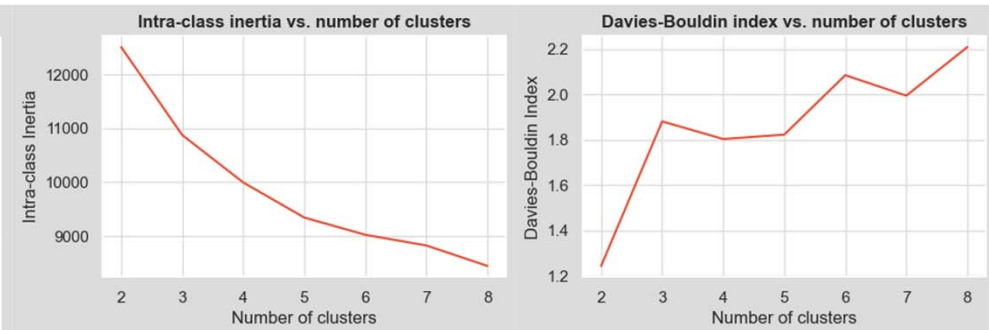
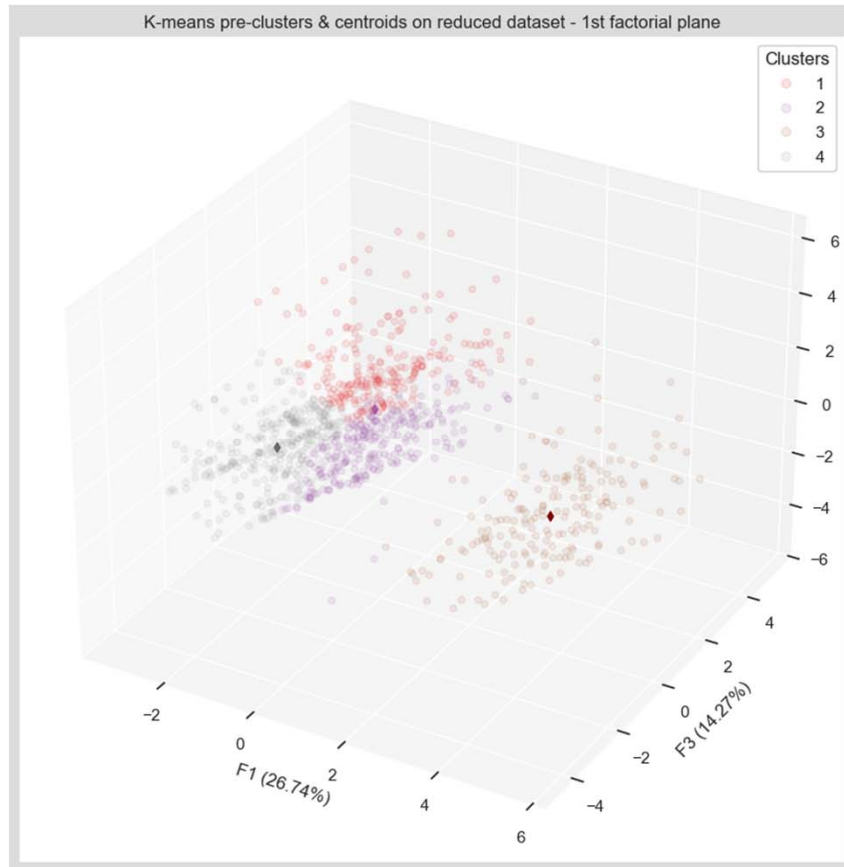


Silhouette Score



2 – MODÉLISATION

2.2 – CRITÈRES MÉTIER & CHOIX DU MODÈLE FINAL – K-MEANS



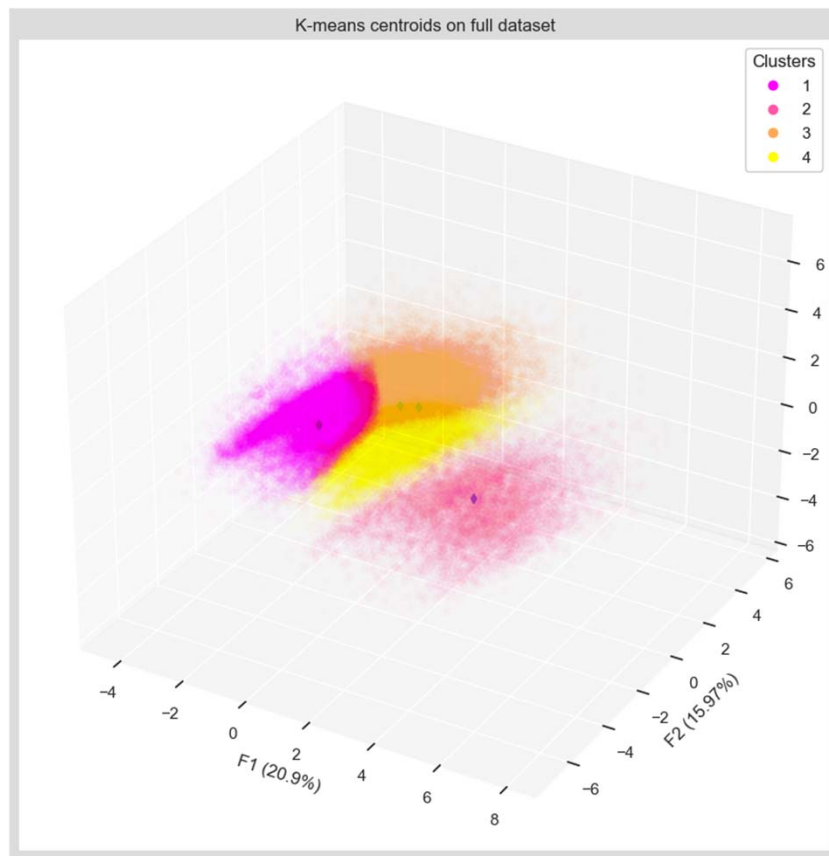
Nb of customer pre-clusters in 1st cluster : 226 . Population : 19771
Nb of customer pre-clusters in 2nd cluster : 279 . Population : 27386
Nb of customer pre-clusters in 3rd cluster : 189 . Population : 12481
Nb of customer pre-clusters in 4th cluster : 306 . Population : 35702
Total customers : 95340 . Total pre-clusters : 1000

Meilleur modèle - compromis :

- Nombre adéquat de clusters pour une stratégie marketing différenciée
- Facilement lisible & interprétable par le métier

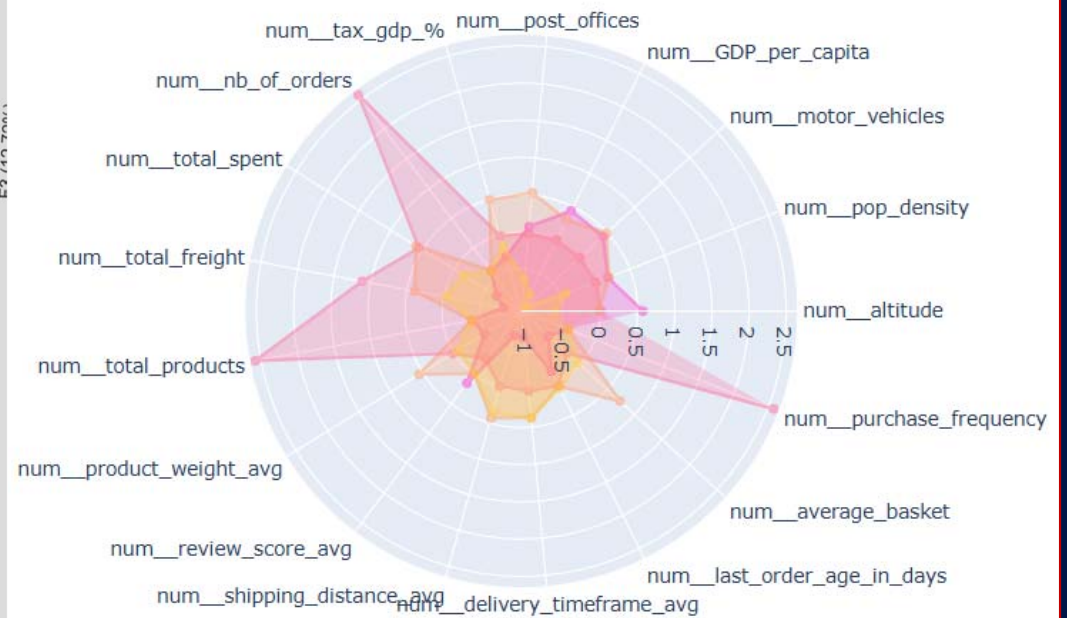
2 – MODÉLISATION

2.2 – CRITÈRES MÉTIER & CHOIX DU MODÈLE FINAL – K-MEANS (SANS PRÉ-CLUSTERING)



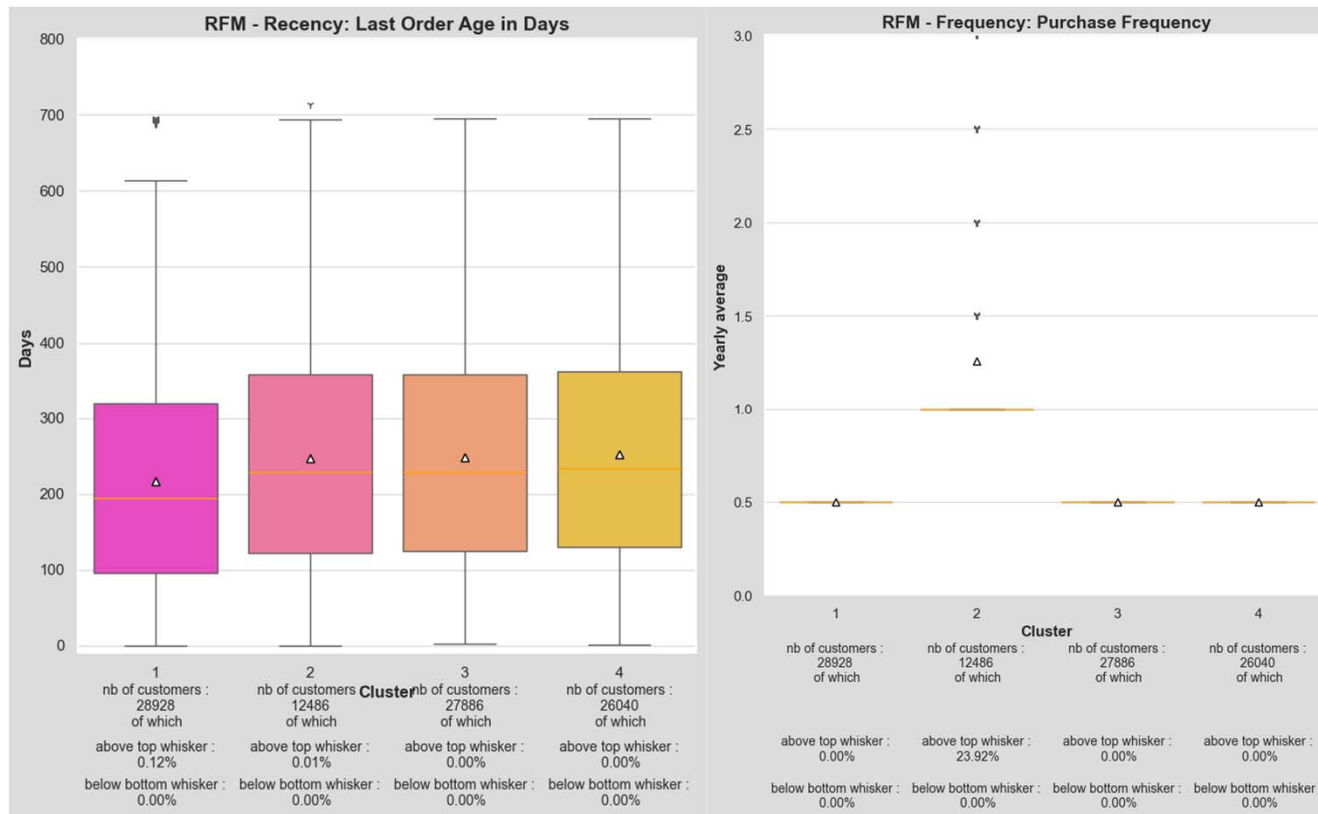
Nb of customers in 1st cluster : 28928
Nb of customers in 2nd cluster : 12486
Nb of customers in 3rd cluster : 27886
Nb of customers in 4th cluster : 26040
Total customers : 95340

Cluster 1 centroid
Cluster 2 centroid
Cluster 3 centroid
Cluster 4 centroid



2 – MODÉLISATION

2.3 – CARACTÉRISTIQUES RFM* DES CLUSTERS & STRATÉGIE MARKETING – RÉCENCE & FRÉQUENCE

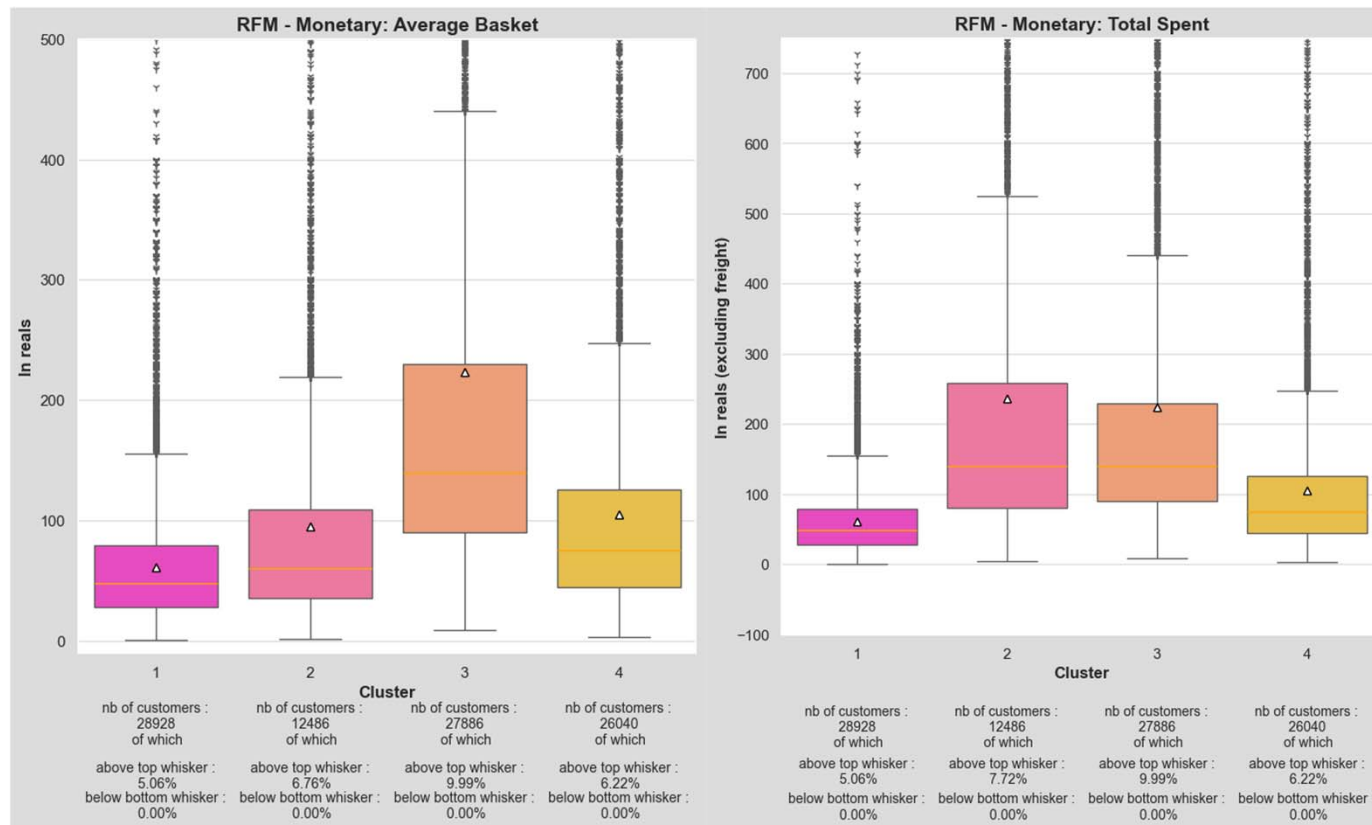


Tests de Kruskal-Wallis & Dunn => différences statistiquement significatives sauf:

- Récence : clusters 2 et 3
- Fréquence : seul le cluster 2 diffère significativement

2 – MODÉLISATION

2.3 – CARACTÉRISTIQUES RFM* DES CLUSTERS & STRATÉGIE MARKETING – MONTANT 1/2

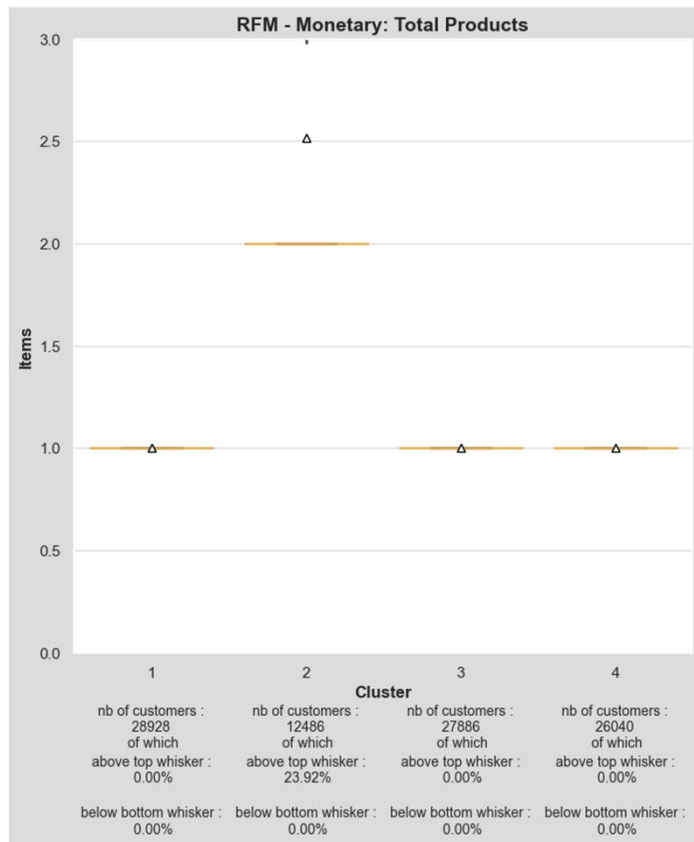


Tests de Kruskal-Wallis & Dunn => différences statistiquement significatives entre tous les clusters sur:

- Panier moyen
- Dépense totale

2 – MODÉLISATION

2.3 – CARACTÉRISTIQUES RFM* DES CLUSTERS & STRATÉGIE MARKETING – MONTANT 2/2



Tests de Kruskal-Wallis & Dunn
=> différences non
statistiquement significatives
entre tous les clusters sauf
**Cluster 2 pour le nombre total
de produits achetés**

2 – MODÉLISATION

2.3 – CARACTÉRISTIQUES RFM* DES CLUSTERS & STRATÉGIE MARKETING – SYNTHÈSE*

Cluster 1 - chenilles



- Clients peu engagés : transactions plus récentes mais peu fréquentes et de faible valeur monétaire
- Stratégie : offres d'entrée de gamme
- Priorité : basse

Cluster 2 - abeilles



- Clients très engagés : transactions moins récentes mais plus nombreuses et plus élevées en montant
- Stratégie : offres personnalisées
- Priorité : haute

Cluster 3 - coccinelles



- Clients sélectifs : transactions moins récentes mais panier moyen et dépense totale les plus élevés
- Stratégie : réduction sur achats en volume ou répétés
- Priorité : moyenne

Cluster 4 – papillons



- Clients testeurs : transactions plus anciennes et de valeur moyenne
- Stratégie : programme de fidélisation
- Priorité : moyenne

3 – EVALUATION DU MODÈLE FINAL

3.1 – STABILITÉ DES CLUSTERS À L'INITIALISATION

PARAMÈTRES:

- 20 itérations / modèle
- random_state = indice itération
- Calcul ARI / itération => moyenne

RÉSULTATS & INTERPRÉTATION:

❖ MINIBATCH K-MEANS

- $ARI_{MOY} = 0.4788$
- Variabilité accrue – échantillons aléatoires & compromis vitesse/précision (taille des échantillons)

❖ K-MEANS

- $ARI_{MOY} = 0.8360$
- Bonne stabilité à l'initialisation & concordance des clusters entre 2 itérations

3 – ÉVALUATION DU MODÈLE FINAL

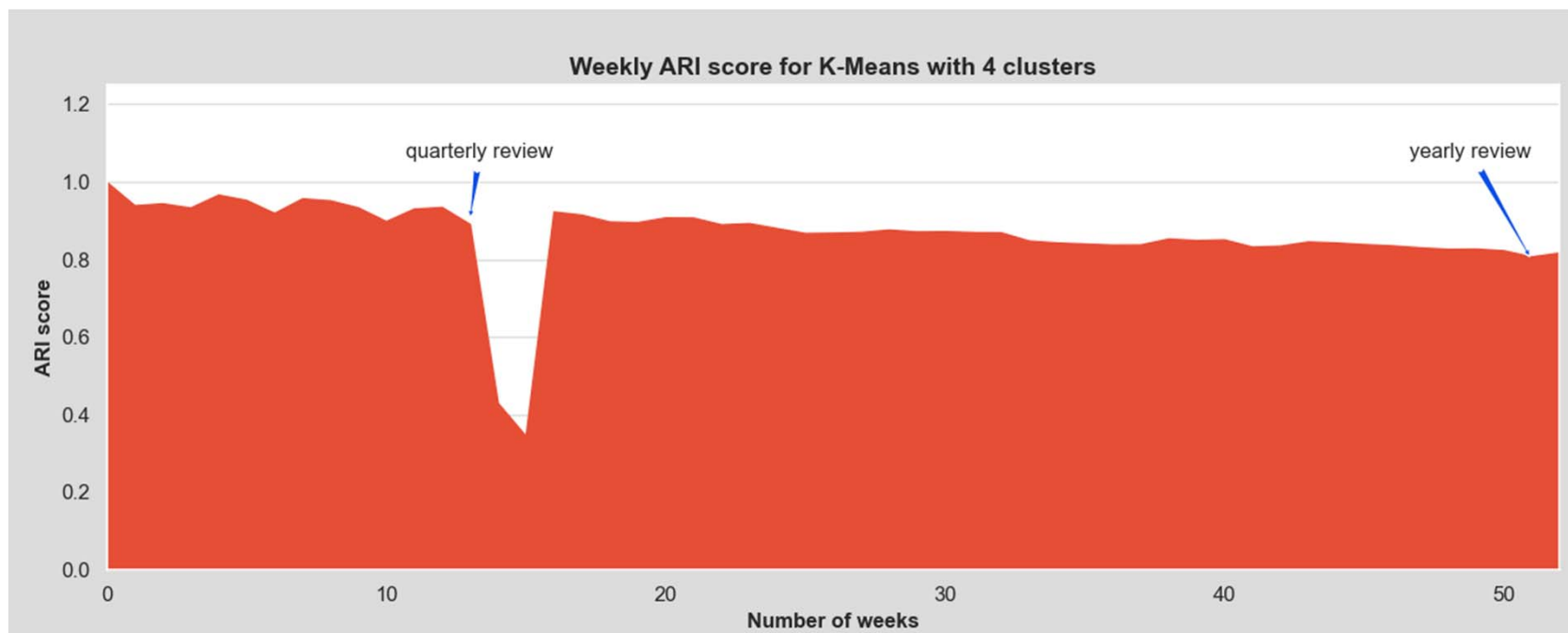
3.2 – CONTRAT DE MAINTENANCE – MÉTHODOLOGIE

MODÈLE => K-MEANS (K=4)

- ❖ T_0 = date dernière transaction client placée (29/08/2018)
 - Entraînement du modèle M_0 à T_0 sur les données D_0
- ❖ T_1, \dots, T_{52} = dates 1, ... , 52 semaines avant T_0
 - Génération d'un fichier de données D_i contenant toutes les transactions depuis le début de l'historique disponible (15/09/2016) jusqu'à D_i
 - Prédiction des clusters C_{0i} du modèle M_0 à T_i sur les données D_i
 - Réentraînement du modèle M_i à T_i sur les données D_i et prédiction des clusters C_{ii}
 - Calcul de l'ARI entre C_{0i} & C_{ii}

3 – ÉVALUATION DU MODÈLE FINAL

3.2 – CONTRAT DE MAINTENANCE – RECOMMANDATIONS



Conclusion & perspectives

❖ Limites de la méthodologie RFM

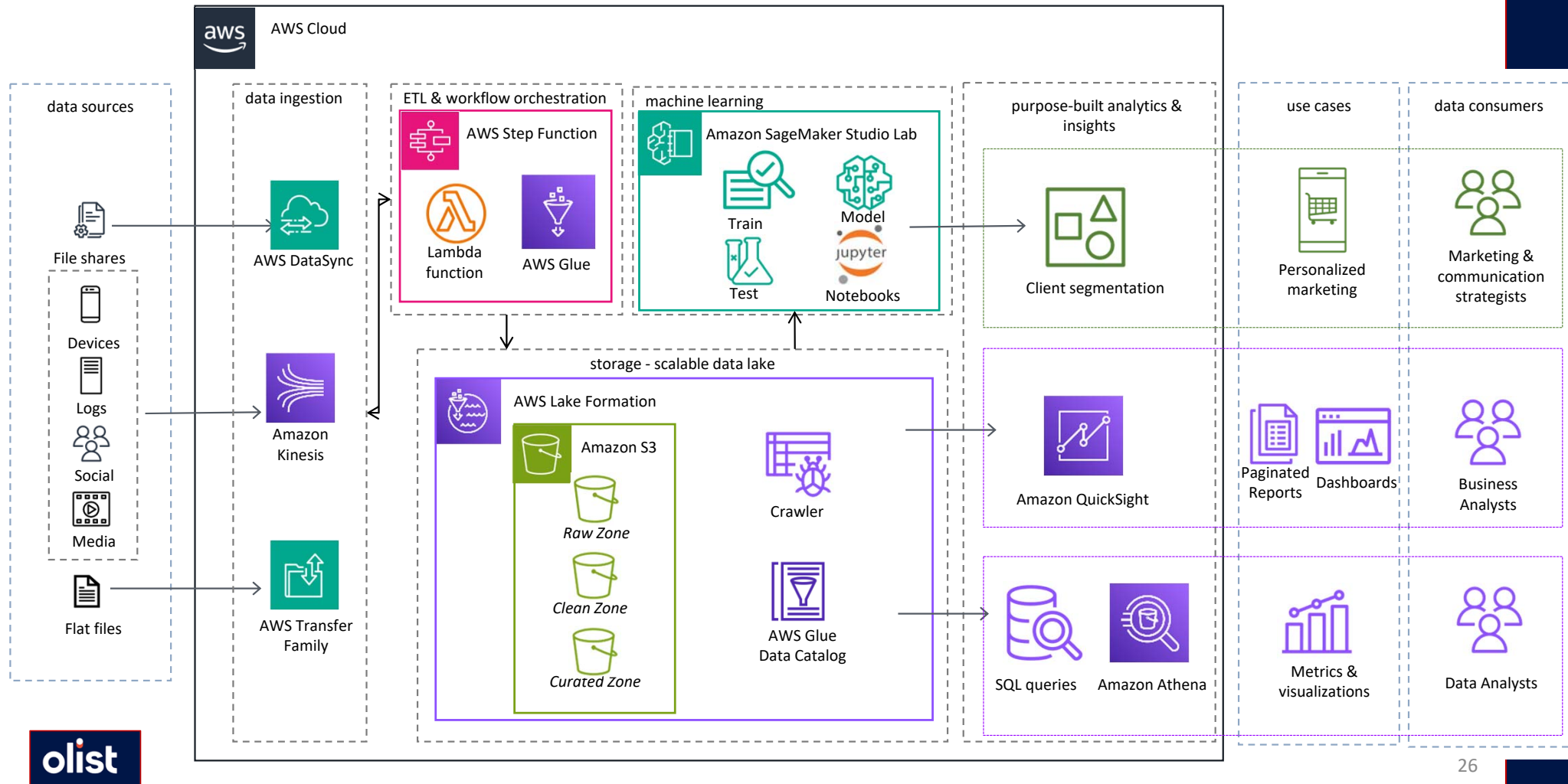
- Récence & fréquence inadaptées à l'analyse de transactions uniques
- Ne mesure pas le « churn risk » ni la loyauté
 - à la “Marketplace” ou à la marque/produit acheté?
- Image incomplète des comportements d'achat qui ignore:
 - la saisonnalité
 - les facteurs psychologiques
 - les influences socio-économico-démographiques

❖ Affiner la stratégie

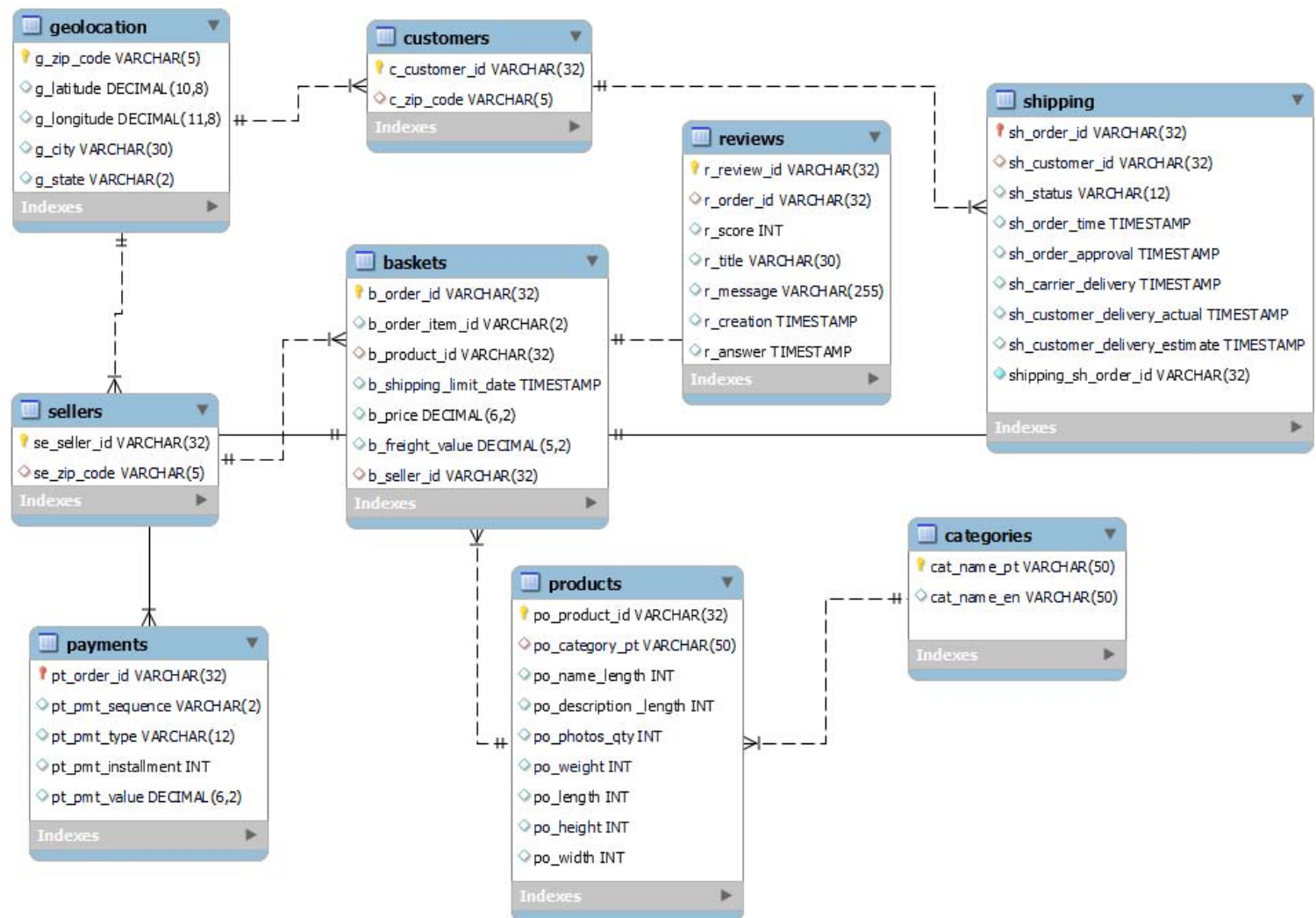
- Identifier les clients potentiels & intégrer les coûts d'acquisition
- Enrichissement des données (démographiques, économiques etc...)*
- Attribution par canal & « marketing-mix »

ANNEXES

ANNEXE 1 – ARCHITECTURE AWS DONNÉES, MODÉLISATION & REPORTING



ANNEXE 2 – ERD



ANNEXE 3 – REQUÊTES SQL / KPI DASHBOARD



3.2%

des commandes sur les 3 derniers mois
ont été livrées au client avec au moins

3 jours de retard



0.5%

des vendeurs ont généré un CA

> 100,000 réals

sur les commandes livrées

Sur les 1,495 nouveaux vendeurs dans les 3
derniers mois, seuls



2

ont vendu

plus de 30 produits

Les 5 codes postaux avec plus de 30 reviews
enregistrant le

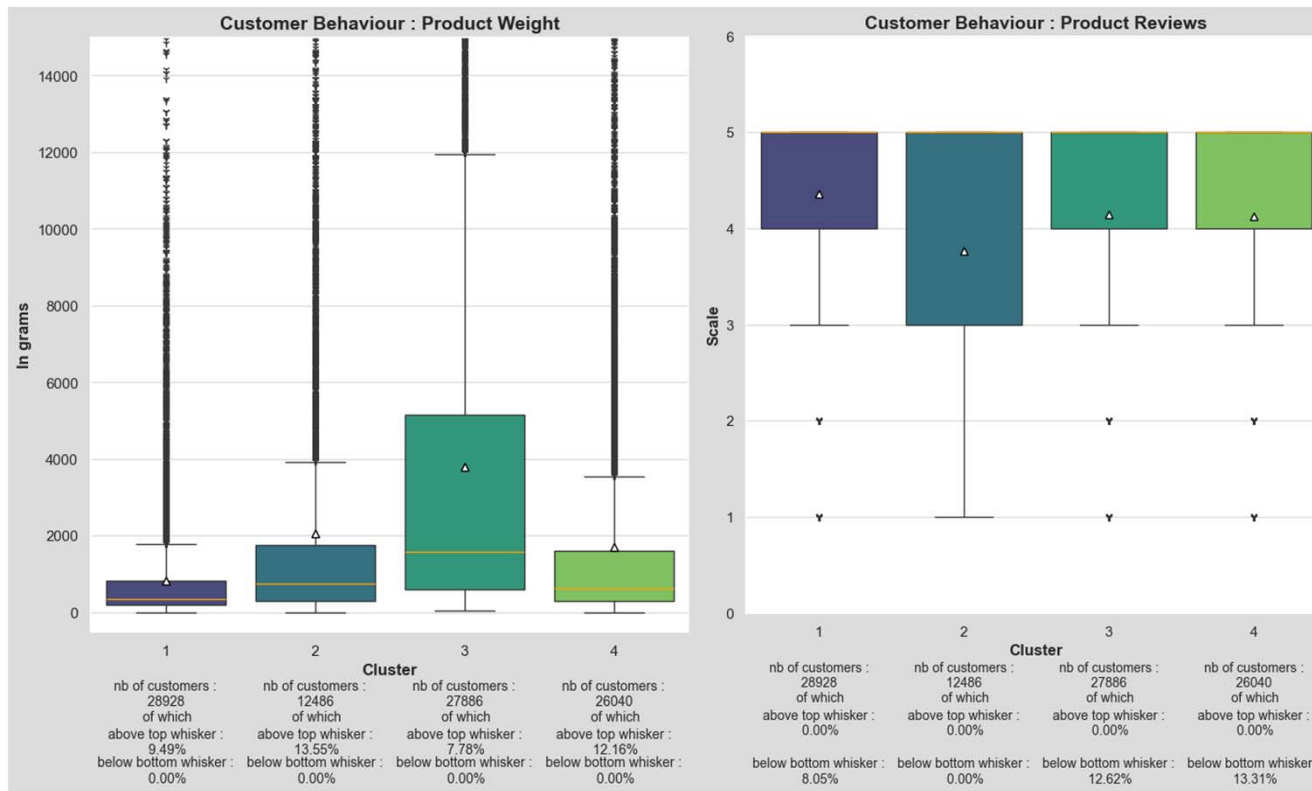


pire score moyen

sur les 12 derniers mois sont tous localisés à

Rio de Janeiro

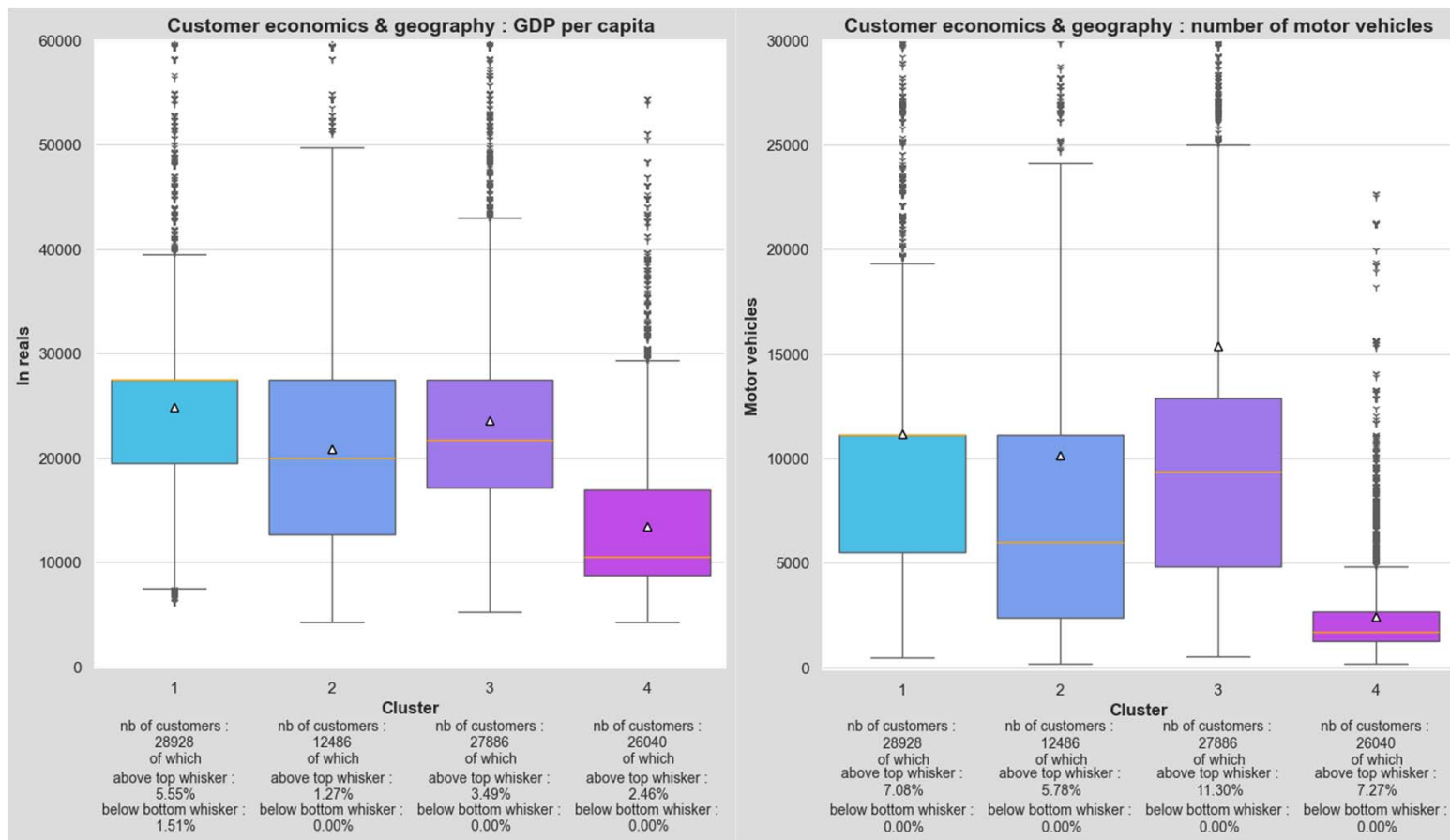
ANNEXE 4A – AUTRES INDICATEURS – COMPORTEMENT CLIENTS



Tests de Kruskal-Wallis & Dunn => différences statistiquement significatives entre tous les clusters sauf:

➤ Product reviews:
Clusters 3 et 4

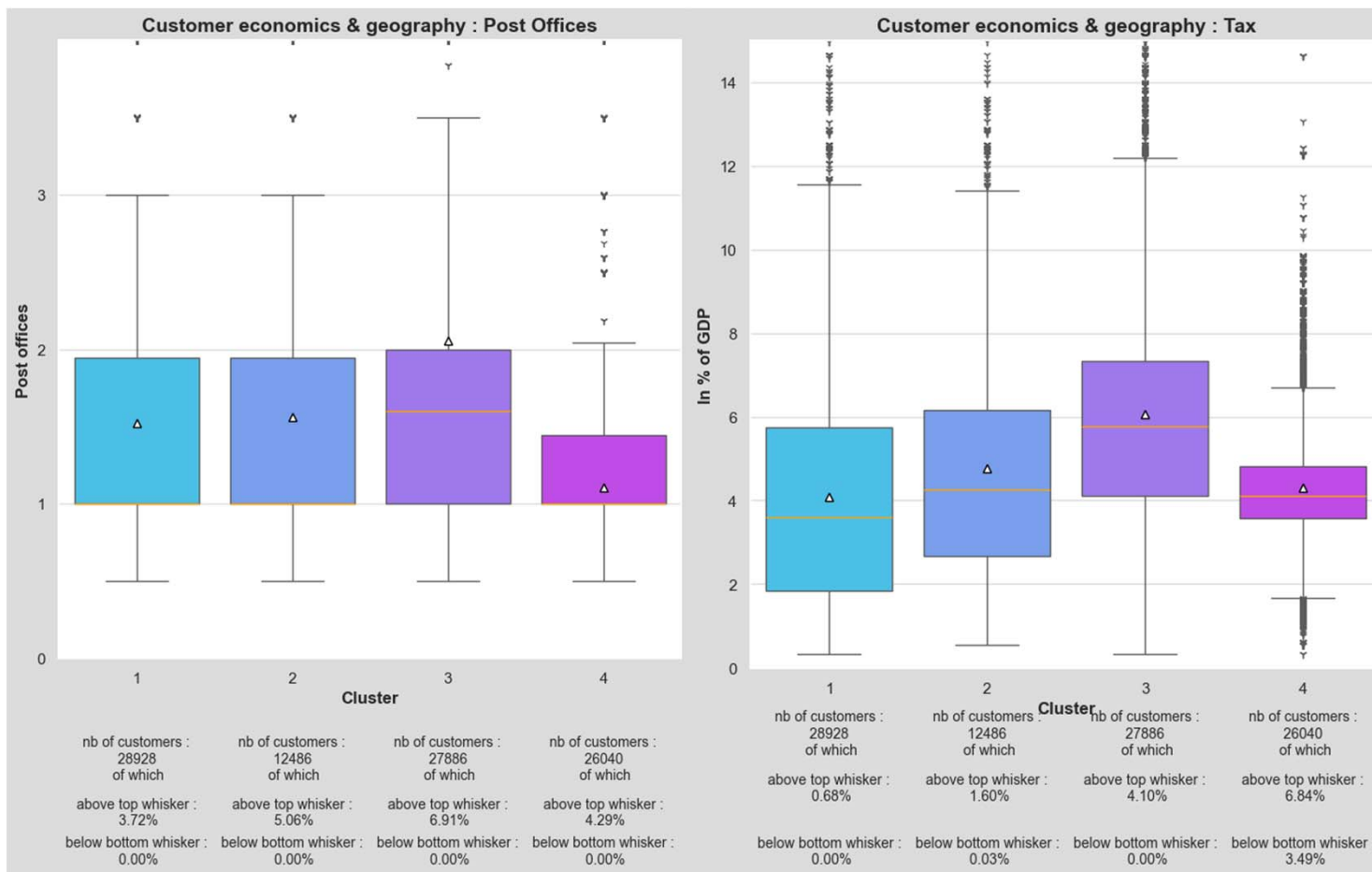
ANNEXE 4B – AUTRES INDICATEURS – ÉCONOMIE & GÉOGRAPHIE*



Tests de Kruskal-Wallis & Dunn => différences statistiquement significatives entre tous les clusters sauf:

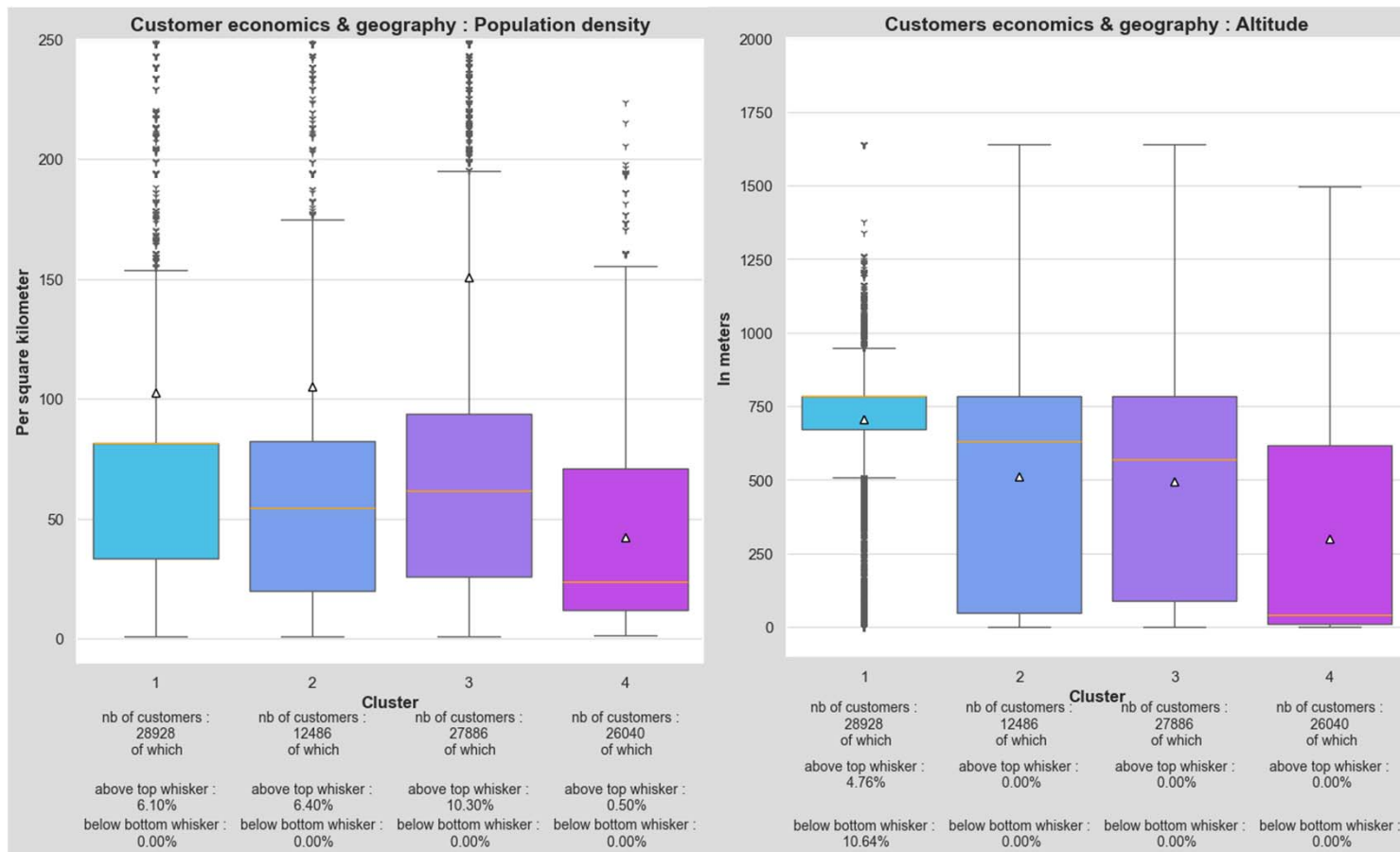
➤ Motor vehicles : Clusters 1 et 3

ANNEXE 4B – AUTRES INDICATEURS – ÉCONOMIE & GÉOGRAPHIE*



Tests de Kruskal-Wallis & Dunn => différences statistiquement significatives entre tous les clusters.

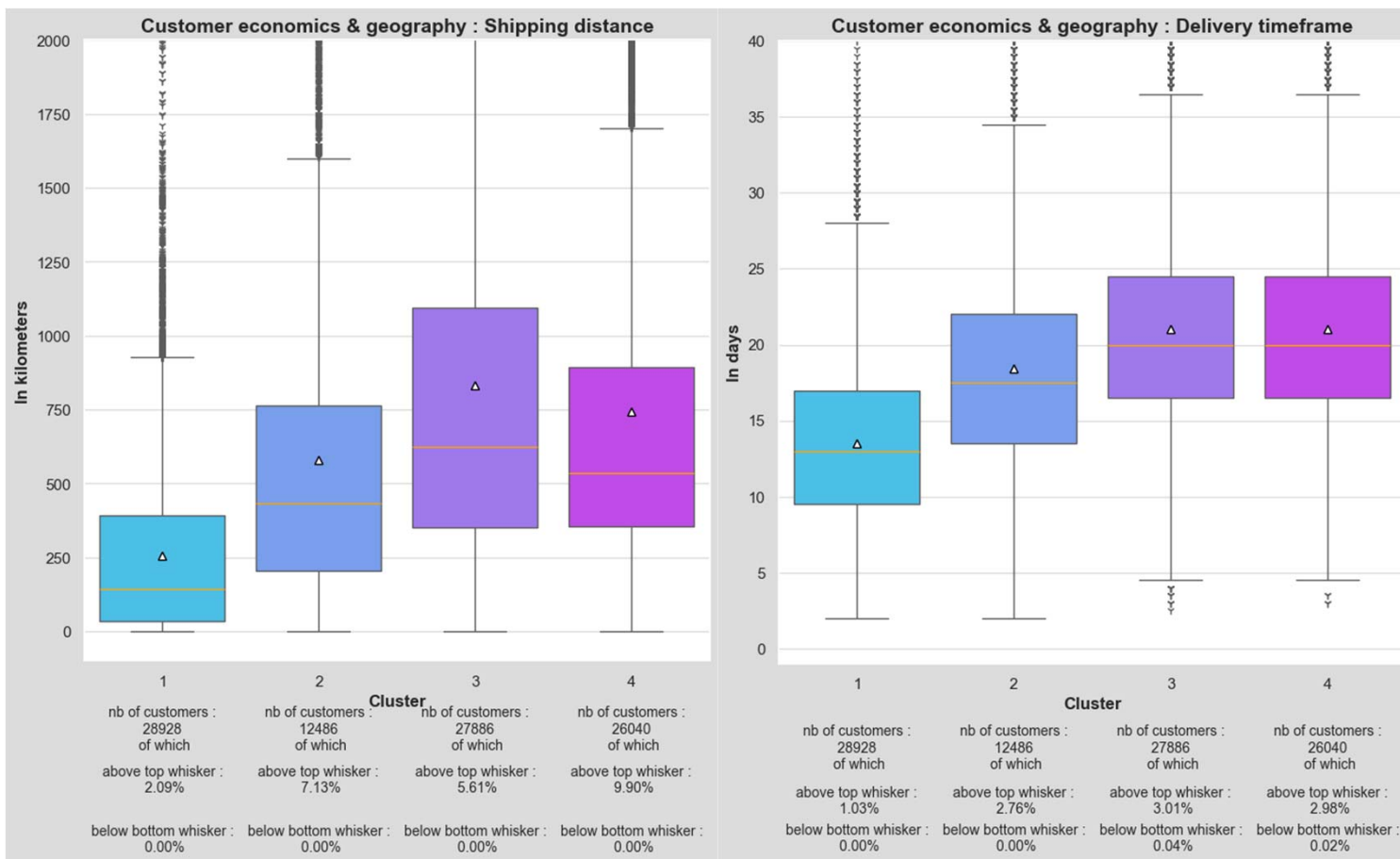
ANNEXE 4B – AUTRES INDICATEURS – ÉCONOMIE & GÉOGRAPHIE*



Tests de Kruskal-Wallis & Dunn => différences statistiquement significatives entre tous les clusters sauf:

➤ Product reviews:
Clusters 3 et 4

ANNEXE 4B – AUTRES INDICATEURS – ÉCONOMIE & GÉOGRAPHIE*



Tests de Kruskal-Wallis & Dunn => différences statistiquement significatives entre tous les clusters sauf:

➤ Delivery timeframe:
Clusters 3 et 4

ANNEXE 4C – SEGMENTATION – CRITÈRES ADDITIONNELS - SYNTHÈSE



Cluster 1 – en villégiature

CSP++, zones à forte densité de population en altitude, délais de livraison plus longs malgré des distances plus courtes, accès moindre aux infrastructures mais forte densité de véhicules à moteur



Cluster 2 - citadins

CSP+, davantage de véhicules et meilleurs accès aux infrastructures que les banlieusards



Cluster 3 – banlieusards

CSP+, bon accès aux infrastructures, gros achats lourds / volumineux, délais de livraison longs, review moyenne la plus basse



Cluster 4 – péri-urbains & ruraux

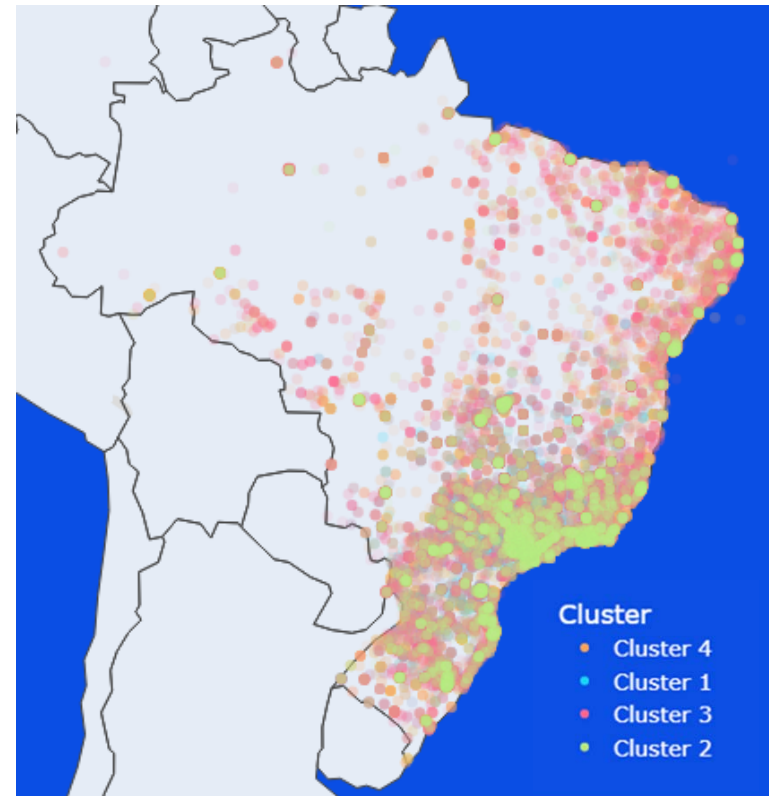
Zones à faible densité de population et basse altitude, peu d'accès aux infrastructures, CSP--, délais d'acheminement longs, grandes distances de livraison

ANNEXE 5 – CUSTOMERS LOCATION

MiniBatch K-Means K=1,000:



K-Means K=4 :



ANNEXE 6 – ANALYSE EN COMPOSANTES PRINCIPALES

